# Streamlining rule-based NLP development and integration with production environment

**HEALTH**
UNIVERSITY OF UTAH

Jianlin Shi, MD, PhD

# A background story

- 1<sup>st</sup> operational NLP project for billing assistant
- Collaborate with EDW, clinical analysts

# Use case

- **Encephalopathy**: A broad term for any brain disease that alters brain function or structure.

- **Missing code** means **revenue loss**

- **Old way:** database keywords search + manual review

# Old Way

(encephalopathy, altered mental, delirium, confusion, AMS, confused, confussed, encephalopathic) ~near((no, encephalopathy), 5, TRUE)  ~near(("not", encephalopathy), 8, TRUE) ~near((denies, encephalopathy), 5, TRUE) ~near((negative, encephalopathy), 5, TRUE) ~near((rule out, encephalopathy), 3, TRUE) ~near((mild, encephalopathy), 3, TRUE) ~near((borderline, encephalopathy), 5, TRUE)...

4

# A web-based interface for code review (Warthog)

| INFERENCE REPORT FOR - *AREGO* - Enceph - 3 (PAST 30 DAYS) | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PAT ID | VISIT NO | NAME | HAR | UNIT | INFERENCE ID | INFERENCE | FIRST INFER DTM | CURRENT INFERENCE STATUS | CURRENT NOTE | STATUS/NOTE HISTORY | |
| ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | | | Enceph - 3 - | 08/03/2017 11:32 | | | ▓▓▓ | ⬇⬆ |

| PAT ID | VISIT NO | SOURCE | RPT TYPE | RPT DESC | TEXT DATE | PREVIEW(RETURNS FIRST FOUND HIGHLIGHTED TERM) | RPT ID |
|---|---|---|---|---|---|---|---|
| ▓▓▓ | ▓▓▓ | EPIC | | | | -minimize nighttime interruptions, at high risk of delirium with ICU hospitalization | ▓▓▓ |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | | Enceph - 3 - | 08/03/2017 11:31 | | | ▓▓▓ | ⬇⬆ |
| ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | | Enceph - 3 - | 08/03/2017 11:31 | | | ▓▓▓ | ⬇⬆ |
| ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | Enceph - 3 - | 08/03/2017 11:31 | | | ▓▓▓ | ⬇⬆ |
| ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | ▓▓▓ | Enceph - 3 - | 08/03/2017 11:32 | | | ▓▓▓ | ⬇⬆ |

# Questions were asked:

- Where to host the program ?
- How to add new inputs?
- How to trigger the execution?
- How to retrieve the NLP outputs?
- How much computing power is needed?
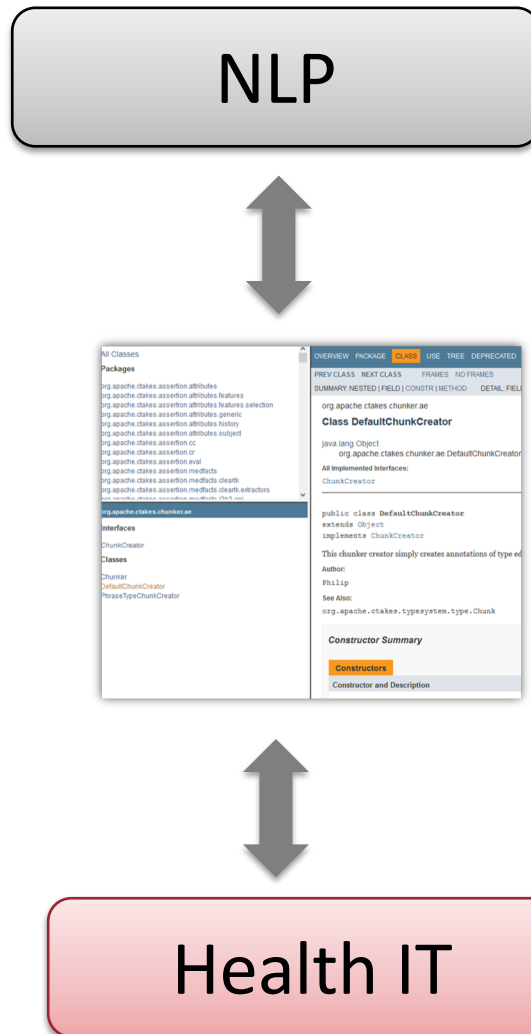- Do we need redo everything for another project?

  .

  .

  .

**HEALTH**
UNIVERSITY OF UTAH

# Questions were asked:

- Where to host the program ?
- How to add new inputs?
- How to trigger the execution?
- How to retrieve the NLP outputs?
- How much computing power is needed?
- <u>Do we need redo everything for another project</u>?
  - .
  - .
  - .

# Research NLP VS Operational NLP

- Reusability — Development cost
- Computing cost
- Interoperability — with EHR and others
- Maintenance
- User support

# One for "all" solution regarding

- Interaction with Health IT systems
- Supporting different NLP tasks—with extensibility

# APIs



- Powerful to access low level functionalities

- Language constrains

- Expert only

- Maintenance strait

# Network Protocols



- Language independent

- Additional storage still needed

- Need client-side implementations

- Maintenance cost

# Database



- Naturally integrated with EHR, EDW

- No technical barrier for EDW staff

- No programming workload for EDW staff

- No additional hurdle for current EDW users

# Database



- Need a unified table structure / schema

# Categories of Clinical NLP Tasks' Output

| Output Categories | Description |
| --- | --- |
| **Mention level** | Identify the concepts or statements in a document |
| **Document level** | Classify a document based on the given conditions |
| **Encounter level** | Classify an encounter based on the given conditions |
| **Patient level** | Classify a patient based on the given conditions |

# Harmonized Database Schema for NLP output

| Output Categories | Column Names | Value Type | Brief description |
|---|---|---|---|
| Mention level | MENTION_RESULT_ID | INTEGER | Mention level results ID |
| | DOCUMENT_RESULT_ID | INTEGER | To be matched to the associated document level results ID |
| | MEN | | |
| | MEN | | |

| | Document level | DOCUMENT_RESULT_ID | INTEGER | Document level result ID |
|---|---|---|---|---|
| | | RUN_ID | INTEGER | Id to identify each NLP executio |
| | | NLP_PIPELIN_ID | INTEGER | Id to identify each NLP pipeline |
| | | DOCUMENT_ID | | |
| | | | | |

| | Encounter level | ENCOUNTER_RESULT_ID | INTEGER | Encounter level result ID |
|---|---|---|---|---|
| | | NLP_INPUT_ID | INTEGER | Id to identify a set of input docu defined by the users, it can be al of the documents in a visit. |
| | | ENCOUNTER_ID | INTEGER | The Id of an encounter |
| | | ENCOUNTER_TYPE | TEXT | The encounter level conclusion |

# Harmonized Database Schema for NLP input

## NLP input cohort table

| Column Name | Value type | Brief Description |
|---|---|---|
| NLP_INPUT_ID | Integer | An ID for a group of documents that needs to be processed |
| PIPELINE_ID | Integer | NLP pipeline ID |
| ..... | | |

| Column Name | Value type | Brief Description |
|---|---|---|
| NLP_INPUT_ID | Integer | An ID for a group of documents that needs to be processed |
| NOTE_ID | Integer | Note ID |
| ..... | | |

## NLP input document table

HEALTH
UNIVERSITY OF UTAH

©UNIVERSITY OF UTAH HEALTH, 2019

# Overall setup

# Workflow of a Generic Rule-based NLP Pipeline



**Patient Level**

**Encounter Level**

**Mention Level**

**Document Level**

Section detector

Sentence segmenter

Named entity recognizer

Context detector

Temporal inferencer

Patient inferencer

Encounter inferencer

Document inferencer

Feature merger

Feature inferencer

18

# One for "all" rule-based NLP

Pipeline 1   Pipeline 2

NLP pipelines

| PipelineID | Component | Rules |
|---|---|---|
| 1 | SectionDetector | |
| 1 | NER | |
| ... | | |
| 2 | SectionDetector | |
| 2 | NER | |
| ... | | |

HEALTH
UNIVERSITY OF UTAH

# From Dev To Production

Pipeline 1   …                                                        Dev Pipeline 1

NLP pipelines

# Overall setup



Rule Configuration

| PipelineID | Component | Rules |
|---|---|---|
| 1 | SectionDetector | |
| 1 | NER | |
| … | | |

# Rule development needs

- Focus on rule development rather than coding
- Operate with different databases
- Support debugging rules
- …

# Tool for Rule Development

# Import and Execute pipeline

# Debug

# Export

| Tasks | Parameter | Value | Description |
|---|---|---|---|
| import | format | | |
| easycie | ehost | data/output/ehost | the directory to save the exported ehost files |
| debug | brat | data/output/brat | the directory to save the exported brat files |
| compare | uima | data/output/xmi | the directory to save the exported uima xmi files |
| export | exportTypes | Concept,Doc_Base | If specified, then only these types will be displayed(separate ... |
| settings | excel | | |
| | sql | | |
| | directory | data/output/excel | The dirctory to save the exported excel files |
| | sampleSize | 600 | If >0, then easyCIE will randomly sample the defined number... |

EasyCIE(__tt_config.xml__)

File    Help

▼ Execute functions

ExportEhost    OpenEhost

**Status:**    Processing Complete.

Reset    Save    Cancel

26

# Streamlining rule-based NLP

- Create a SQL query

**Define Cohort**

**Develop Rule-Based NLP**

- Only need to care about rules

- Turn on a switch

**Publish Pipeline**

# Thank you!

# n-Trie demo

Element 1          Element 2          Element 3

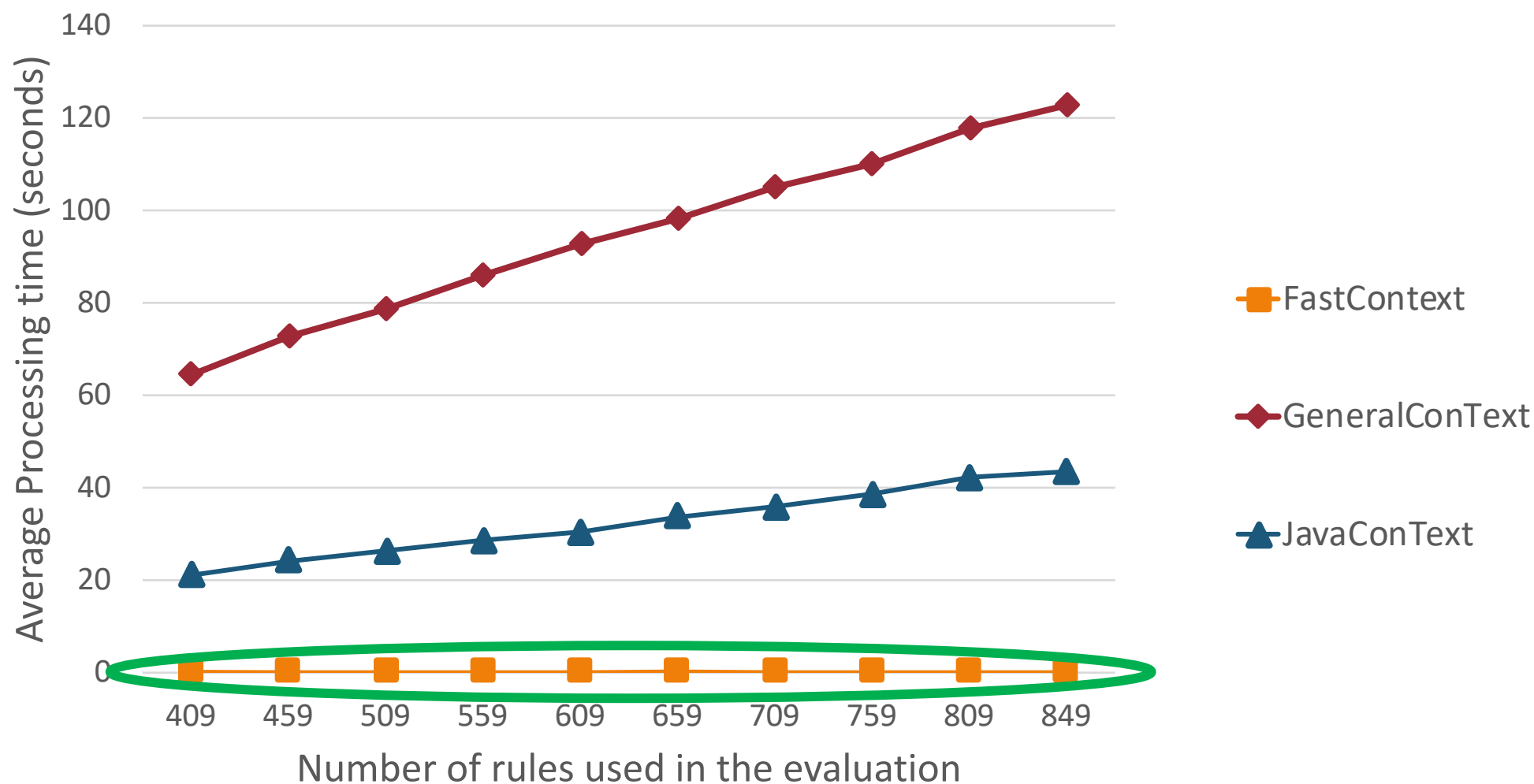| rule out | rule → out |
| without | without |
| can't rule out | can't → rule → out |
| can't exclude | can't → exclude |

# FastContext — Speed Evaluation



Average processing time of the whole dataset in 200 runs

30

# Encephalopathy pipeline evaluation

Totally sampled 665 visits,
including   8068 documents

|  | Found encephalopathy mentions (# of visit) | Verified by review (# of visits) |
|---|---|---|
| Old Warthog | 50 | 13 |
| NLP powered Warthog | 208 | 178 |

Precision improves 231%,
Recall improves 1269%

31

HEALTH
UNIVERSITY OF UTAH

# Methods---PE Identifier Evaluation

- A **local** dataset (400 annotated CT pulmonary angiography reports)
- **Stanford** dataset (944 reports)
- **PEfinder** dataset (859 reports)
- Compared with
  - Intelligent Word Embedding (IWE)[1]
  - PEfinder[2]

1. Banerjee I, Chen MC, Lungren MP, et al. Radiology report annotation using intelligent word embeddings: Applied to multi-institutional chest CT cohort. J Biomed Inform 2018;77:11–20. doi:10.1016/j.jbi.2017.11.012
2. Chapman BE, Lee S, Kang HP, et al. Document-level classification of CT pulmonary angiography reports based on an extension of the ConText algorithm. J Biomed Inform 2011;44:728–37. doi:10.1016/j.jbi.2011.03.011

32

# PE Identifier Evaluation